# Developing a Born-Digital Preservation Workflow

## Bill Donovan   &   Jack Kearney

Digital Imaging & Curation Manager     Audiovisual Archives Assistant

Thomas P. O'Neill, Jr. Library     John J. Burns Library

Boston College Libraries

July 23, 2014

# Goals

1. Develop a systematic approach to digital preservation (DP) of born-digital collections.

2. Gain experience with a variety of DP hardware and software and figure out practical protocols --- how does all of this stuff work?

3. Use a real-life example, the electronic records of the Mary O'Hara papers. (MOH)

# Mary O'Hara electronic records



An Irish soprano and harpist of international renown, Mary O'Hara has appeared on many of the world's major stages, including Royal Albert Hall, New York's Carnegie Hall, Sydney Opera House, and Toronto's Roy Thompson Hall.

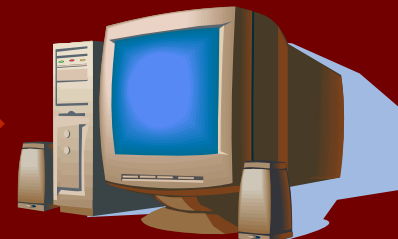Hard drive donated to the Burns Library Irish Music Center at Boston College.
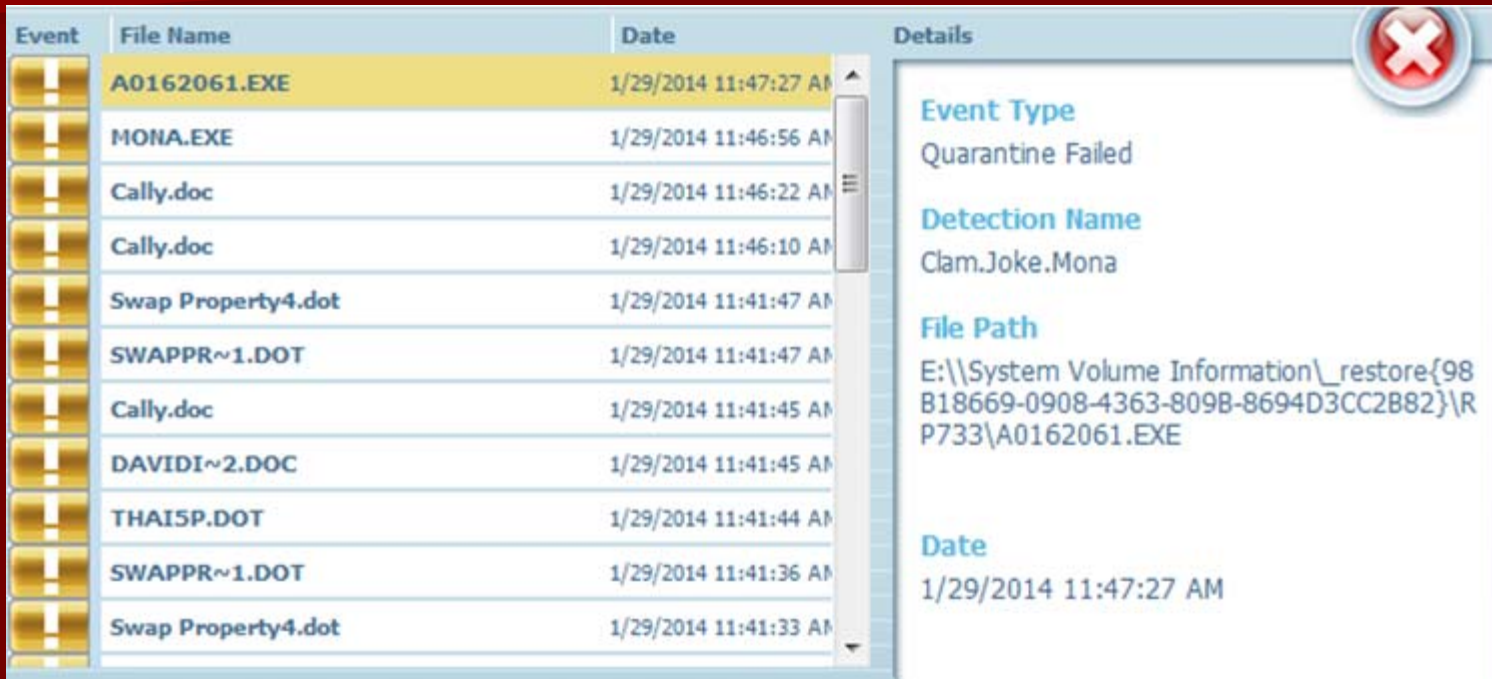
# Protecting the original

Used a forensic "write-blocker." Permits reading but not writing. Prevents changes to the files on the hard drive.

USB

USB

DP Workstation

# But is the original virus-free?

| Event | File Name | Date | Details |
|---|---|---|---|
| | A0162061.EXE | 1/29/2014 11:47:27 AM | **Event Type** |
| | MONA.EXE | 1/29/2014 11:46:56 AM | Quarantine Failed |
| | Cally.doc | 1/29/2014 11:46:22 AM | |
| | Cally.doc | 1/29/2014 11:46:10 AM | **Detection Name** |
| | Swap Property4.dot | 1/29/2014 11:41:47 AM | Clam.Joke.Mona |
| | SWAPPR~1.DOT | 1/29/2014 11:41:47 AM | **File Path** |
| | Cally.doc | 1/29/2014 11:41:45 AM | E:\\System Volume Information\_restore{98 B18669-0908-4363-809B-8694D3CC2B82}\R P733\A0162061.EXE |
| | DAVIDI~2.DOC | 1/29/2014 11:41:45 AM | |
| | THAI5P.DOT | 1/29/2014 11:41:44 AM | |
| | SWAPPR~1.DOT | 1/29/2014 11:41:36 AM | **Date** |
| | Swap Property4.dot | 1/29/2014 11:41:33 AM | 1/29/2014 11:47:27 AM |

Decision: Delete 33 infected files from the hard drive (after first creating a comprehensive inventory of all files on the drive)
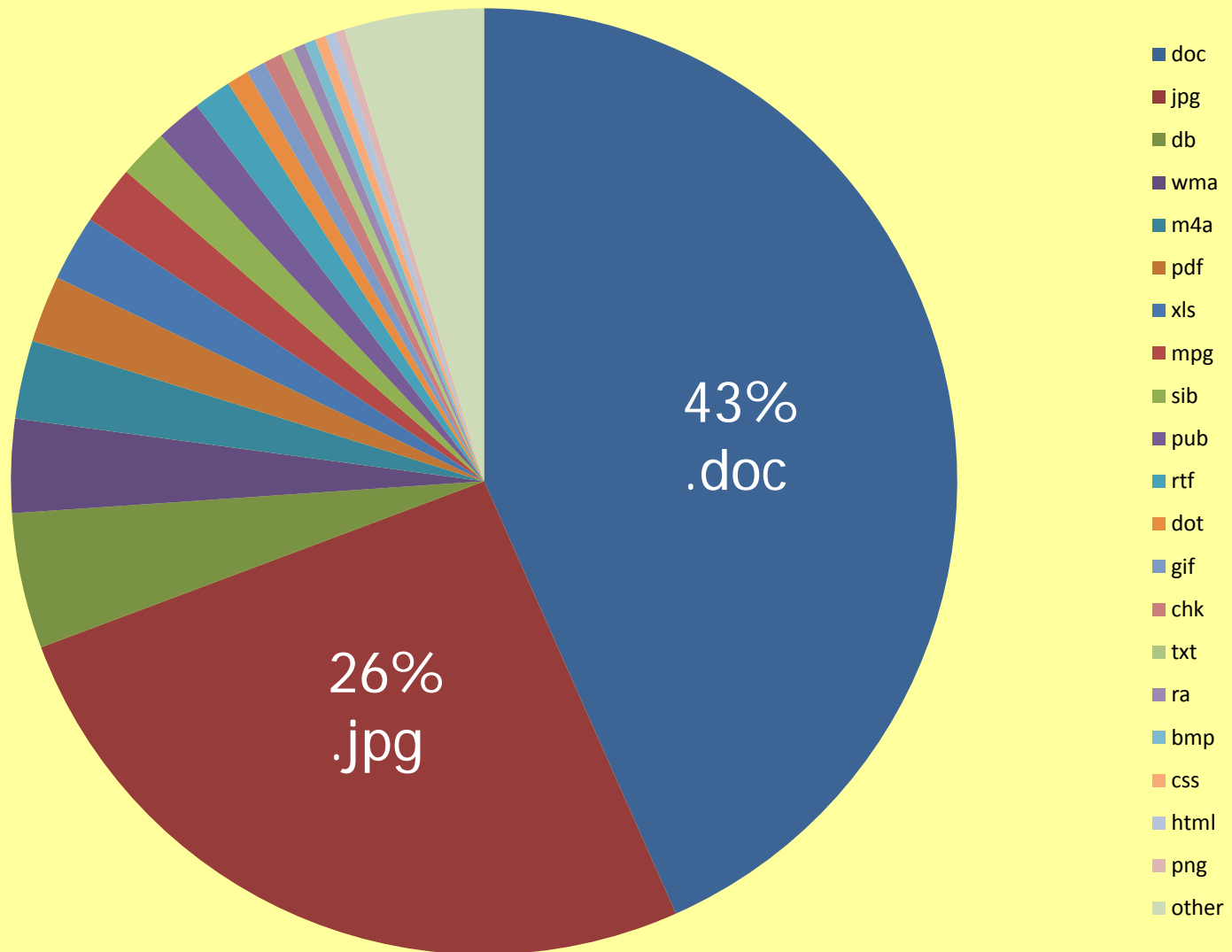
# What's on the MOH hard drive?
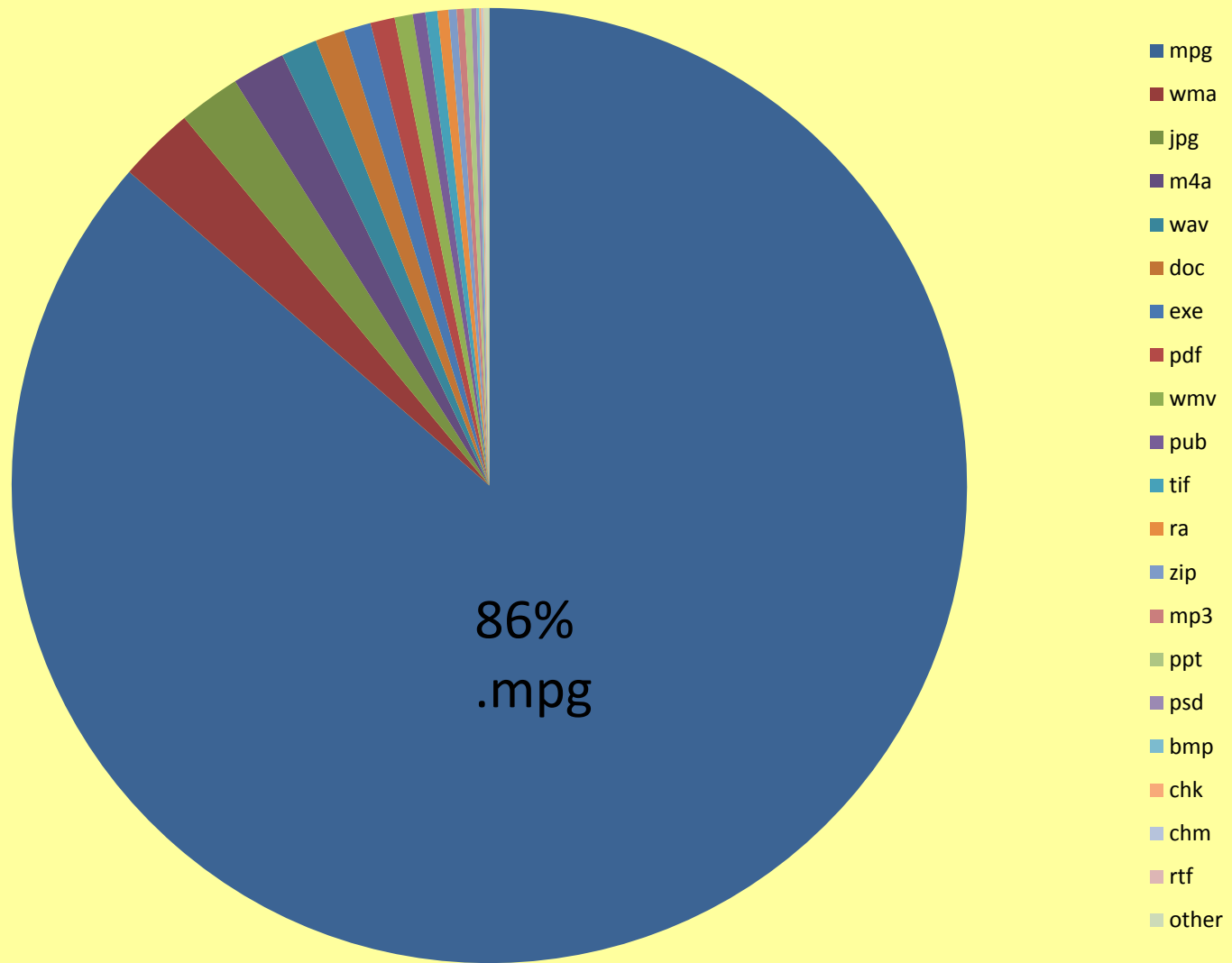
Inventory taken with this Unix command:

find *directory-name* -type f -exec ls -l {} ; >c:\data\MOH\inventory.txt

21,988 files that are on the hard drive, for a total of 104.3 GB

# 20 most frequent file types, plus other



43%
.doc

26%
.jpg

- doc
- jpg
- db
- wma
- m4a
- pdf
- xls
- mpg
- sib
- pub
- rtf
- dot
- gif
- chk
- txt
- ra
- bmp
- css
- html
- png
- other

# 20 file types with most number of bytes, plus other

86%
.mpg

- mpg
- wma
- jpg
- m4a
- wav
- doc
- exe
- pdf
- wmv
- pub
- tif
- ra
- zip
- mp3
- ppt
- psd
- bmp
- chk
- chm
- rtf
- other

# Duplicate files



"8.3" constraint

Decision: keep only the copy with the more complete filename

# De-duping software

# Initial checksums

Decide which checksum algorithm to use

Compute initial checksums for each set of files

Compare checksums to verify original files and copies are <u>identical</u>

| | A | B | C |
|---|---|---|---|
| 1 | external | internal-working | |
| 2 | 0491a43250f36f5cab7e47ffb4f59691 | 0491a43250f36f5cab7e47ffb4f59691 | 0 |
| 3 | 824a2cf8b9e918f66a4ec0c17ab3e4ed | 824a2cf8b9e918f66a4ec0c17ab3e4ed | 0 |
| 4 | c57840c3e45130a86390c5bbd921aeaf | c57840c3e45130a86390c5bbd921aeaf | 0 |
| 5 | e59574a768d30e5d516c601f09cba5d5 | e59574a768d30e5d516c601f09cba5d5 | 0 |
| 6 | 824a2cf8b9e918f66a4ec0c17ab3e4ed | 824a2cf8b9e918f66a4ec0c17ab3e4ed | 0 |
| 7 | c57840c3e45130a86390c5bbd921aeaf | c57840c3e45130a86390c5bbd921aeaf | 0 |
| 8 | e59574a768d30e5d516c601f09cba5d5 | e59574a768d30e5d516c601f09cba5d5 | 0 |
| 9 | fdc1641c56df01318498024599ac09f9 | fdc1641c56df01318498024599ac09f9 | 0 |
| 10 | c354c0f1f4054d25ef28e60ae8394b4d | c354c0f1f4054d25ef28e60ae8394b4d | 0 |
| 11 | eb3a98ae76185ca2933fe9592470a14b | eb3a98ae76185ca2933fe9592470a14b | 0 |

For final-state collections, use automated tool for integrity-checking

# Odd file/folder names

**HDD MEDIA PLAYER**
File folder

**MaryO'HaraDesktop**
File folder

**MO'H_Hard Disc**
File folder

**MO'Hara_MediaPlayer**
File folder

| | | | |
|---|---|---|---|
| # pound | < left angle bracket | $ dollar sign | + plus sign |
| % percent | > right angle bracket | ! exclamation point | ` backtick |
| & ampersand | * asterisk | ' single quotes | \| pipe |
| { left bracket | ? question mark | " double quotes | = equal sign |
| } right bracket | / forward slash | : colon | |
| \ back slash | blank spaces | @ at sign | |

Some "forbidden characters"

http://www.mtu.edu/umc/services/web/cms/characters-avoid/

Had to "escape" these characters in folder name in our Unix commands.

**Decision**: remediate folder and file names, <u>but only for the working copies</u>.

# Renamer: remediating folder/file names

# Any files off-limits or expendable?

- Confidential information
  - Social security numbers
  - Financial information
- Files that have nothing to do with MOH per se
  - System Files
- Files that have no archival value
  - Thumbs.db

# Personally Identifiable Information (PII)



Policy decisions to be made by archivists:
- Based upon the PII findings, which files will eventually made open to the public, or not?

# Any proprietary file formats?

Normalize using Xena ? (see Gregory 2010)



For RealAudio files, normalize with dBpoweramp?

Policy decision:  Which files to normalize, and which formats to preserve?

# File formats --- identify, validate, and extract metadata --- using FITS



http://projects.iq.harvard.edu/fits/fits-processing

# Output of FITS --- identity of the file and the current tool versions in FITS

...

    &lt;identity format="Tagged Image File Format" mimetype="image/tiff" toolname="FITS" toolversion="0.6.2"&gt;

      &lt;tool toolname="Jhove" toolversion="1.5" /&gt;

      &lt;tool toolname="file utility" toolversion="5.03" /&gt;

      &lt;tool toolname="Exiftool" toolversion="9.06" /&gt;

      &lt;tool toolname="Droid" toolversion="3.0" /&gt;

      &lt;tool toolname="NLNZ Metadata Extractor" toolversion="3.4GA" /&gt;

      &lt;tool toolname="ffident" toolversion="0.2" /&gt;

...

# Output of FITS --- file information

...

    &lt;size toolname="Jhove" toolversion="1.5"&gt;1795770&lt;/size&gt;

    &lt;creatingApplicationName toolname="Jhove" toolversion="1.5"&gt;Omniscan 11.12 SR2 Build13&lt;/creatingApplicationName&gt;

    &lt;lastmodified toolname="Exiftool" toolversion="9.06" status="SINGLE_RESULT"&gt;2013:08:14 14:15:38-04:00&lt;/lastmodified&gt;

    &lt;filepath toolname="OIS File Information" toolversion="0.1" status="SINGLE_RESULT"&gt;C:\DATA\FITS_test_folder\fits_test_imagefile.tif&lt;/filepath&gt;

    &lt;filename toolname="OIS File Information" toolversion="0.1" status="SINGLE_RESULT"&gt;C:\DATA\FITS_test_folder\fits_test_imagefile.tif&lt;/filename&gt;

    &lt;md5checksum toolname="OIS File Information" toolversion="0.1" status="SINGLE_RESULT"&gt;ccfca47fb4f2597c04e299c99f4043ce&lt;/md5checksum&gt;

    &lt;fslastmodified toolname="OIS File Information" toolversion="0.1" status="SINGLE_RESULT"&gt;1376504138000&lt;/fslastmodified&gt;

...

# Output of FITS --- file status

```
<filestatus>
    <well-formed toolname="Jhove" toolversion="1.5"
status="SINGLE_RESULT">true</well-formed>

    <valid toolname="Jhove" toolversion="1.5" status="SINGLE_RESULT">true</valid>
  </filestatus>
```

From JHOVE website (http://jhove.sourceforge.net/):

A digital object is well-formed if it meets the purely syntactic requirements for its format.

An object is valid if it is well-formed and it meets additional semantic-level requirements.

# Output of FITS --- metadata

...

&lt;compressionScheme toolname="Jhove" toolversion="1.5"&gt;Uncompressed&lt;/compressionScheme&gt;

&lt;imageWidth toolname="Jhove" toolversion="1.5"&gt;1598&lt;/imageWidth&gt;

&lt;imageHeight toolname="Jhove" toolversion="1.5"&gt;373&lt;/imageHeight&gt;

&lt;colorSpace toolname="Jhove" toolversion="1.5"&gt;RGB&lt;/colorSpace&gt;

&lt;referenceBlackWhite toolname="Jhove" toolversion="1.5" status="SINGLE_RESULT"&gt;0.0 255.0 0.0 255.0 0.0 255.0&lt;/referenceBlackWhite&gt;

&lt;iccProfileName toolname="Exiftool" toolversion="9.06" status="SINGLE_RESULT"&gt;sRGB IEC61966-2.1&lt;/iccProfileName&gt;

&lt;orientation toolname="Jhove" toolversion="1.5"&gt;normal*&lt;/orientation&gt;

&lt;samplingFrequencyUnit toolname="Jhove" toolversion="1.5" status="CONFLICT"&gt;in.&lt;/samplingFrequencyUnit&gt;

&lt;samplingFrequencyUnit toolname="Exiftool" toolversion="9.06" status="CONFLICT"&gt;inches&lt;/samplingFrequencyUnit&gt;

&lt;xSamplingFrequency toolname="Jhove" toolversion="1.5"&gt;300&lt;/xSamplingFrequency&gt;

&lt;ySamplingFrequency toolname="Jhove" toolversion="1.5"&gt;300&lt;/ySamplingFrequency&gt;

&lt;bitsPerSample toolname="Jhove" toolversion="1.5"&gt;8 8 8&lt;/bitsPerSample&gt;

&lt;samplesPerPixel toolname="Jhove" toolversion="1.5"&gt;3&lt;/samplesPerPixel&gt;

&lt;imageProducer toolname="Jhove" toolversion="1.5" status="SINGLE_RESULT"&gt;Zeutschel Omniscan 11&lt;/imageProducer&gt;

&lt;scanningSoftwareName toolname="Jhove" toolversion="1.5"&gt;Omniscan 11.12 SR2 Build13&lt;/scanningSoftwareName&gt;

...

# Workflow: checklist

- ☑ Take inventory
- ☑ Check for viruses
- ☑ Make back-up copies
- ☑ Leave originals as-is
- ☑ Fingerprint each file (checksum)
- ☑ Protect privacy (personally identifiable information)
- ☑ Fix unorthodox folder/file names
- ☑ Avoid proprietary file formats
- ☑ Validate file formats
- ☑ Weed out extraneous/duplicate files
- ☑ Integrate into collection
- ☑ Provide access to users

# Fixity: automated integrity-checking for final state collections



http://www.avpreserve.com/

# Distributed Digital Preservation

LOCKSS-based MetaArchive Cooperative

- At least 6 copies, geographically dispersed
- Fixity-checking (not just back-ups)

# Keep track of DP actions

- File migrations
    - Obsolete file formats
    - Proprietary file formats

- Metadata changes

Possibilities for documenting DP events:

- ArchiveSpace

- PREMIS Event Service

# Future Plans

Additional electronic records from the Mary O'Hara Papers (e.g. data DVDs or CDs)

Replicate DP system but portable

DP in a box

## Beyond preservation...

- Hand off to our archivists (policy questions)

- Provide access to the MOH electronic records

- Create links within the MOH finding Aid

# Additional Resources

ASERL Webinars re: Digital Preservation (Spring 2013)
http://www.aserl.org/intro-dp-2013/

OCLC Research "Demystifying Born Digital" Reports (2012-13)
http://www.oclc.org/research/publications/library/2012/2012-06r.html

BitCurator Project White Paper: "Bitstreams to Heritage:
Putting Digital Forensics into Practice in Collecting
Institutions" (September 2013)
http://www.bitcurator.net/docs/bitstreams-to-heritage.pdf

ARL SPEC Kit 329: Managing Born-Digital Special
Collections and Archival Materials (August 2012)
http://publications.arl.org/Managing-Born-Digital-Special-Collections-
and-Archival-Materials-SPEC-Kit-329

# Q & A

Bill Donovan       bill.donovan@bc.edu

Jack Kearney       kearneyj@bc.edu

# Archival policy questions

- Preserve just the digital files or the entire disk image?

- Delete the virus-infected files?

- Save the duplicate 8.3 files?

- Delete extraneous files? (And, define extraneous)

- Decide fate of PII files:
    - Credit cards
    - Bank accounts
    - Social Security

- Normalize file formats?

- Which files to preserve?
    - External hard drive
    - As-is copies of original files
    - Remediated copies (And, define remediated)

- What sorts of documentation to preserve and how?

# DP Workstation

hardware:

- PC desktop computer with 64-bit Windows 7
- 4 GB RAM and a 465 GB internal hard drive
- UltraKit III + FireWire Write-Blocker

software:

- Cygwin 64 Terminal (md5sum)
- Immunet 3 powered by ClamAV
- Identity Finder
- FITS
- Fixity
- Access Data FTK Imager
- dBpoweramp Music Converter
- HxD Hex Editor

security measures:

- authorized personnel only room
- security cable for workstation
- need-to-know only username/password